

(19)日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開平9-230893

(43)公開日 平成9年(1997)9月5日

(51)Int.Cl. ⁸	識別記号	庁内整理番号	F I	技術表示箇所
G 1 0 L 5/04			G 1 0 L 5/04	D
				F
3/00			3/00	H

審査請求 未請求 請求項の数 5 O L (全 9 頁)

(21)出願番号 特願平8-35291

(22)出願日 平成8年(1996)2月22日

(71)出願人 000102728

エヌ・ティ・ティ・データ通信株式会社
東京都江東区豊洲三丁目3番3号

(72)発明者 林 慶士

東京都江東区豊洲三丁目3番3号 エヌ・
ティ・ティ・データ通信株式会社内

(72)発明者 保理江 高志

東京都江東区豊洲三丁目3番3号 エヌ・
ティ・ティ・データ通信株式会社内

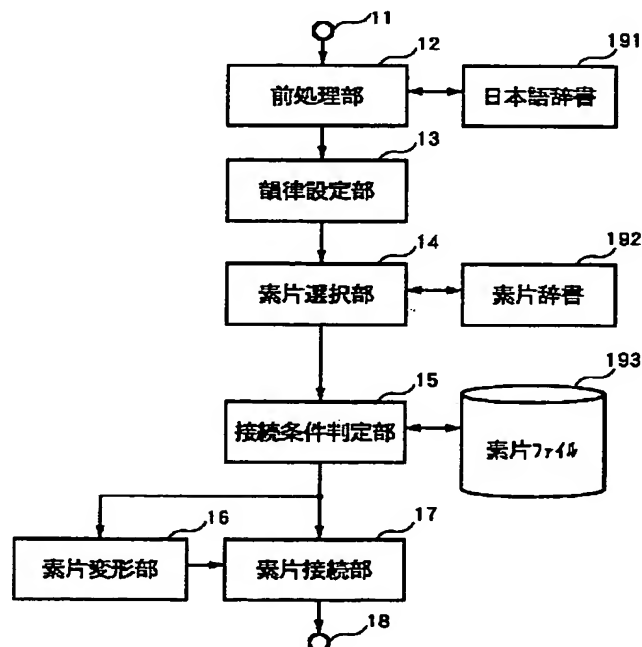
(74)代理人 弁理士 鈴木 正剛

(54)【発明の名称】 規則音声合成方法及び音声合成装置

(57)【要約】

【課題】 合成単位に属する最適素片を選択して接続する音声合成装置において、合成時の音韻環境等の一致度を高めて自然音声に近い音声进行合成する。

【解決手段】 入力文字列を音韻単位に分解する前処理部12と、自然音声から発話単位ごとに切り出された複数の波形素片、及び複数の単音節と各単音節の韻律パラメタを格納する手段182、183と、素片選択部14と、選択された波形素片を入力文字列の順に接続して合成音声进行生成する素片接続手段15、16、17とを含んで音声合成装置を構成する。素片選択部は、音韻環境に対して、その抽出環境及び韻律パラメタの双方を基準として、音韻単位に対応する最適な素片を選択する。



【特許請求の範囲】

【請求項 1】 音韻単位に属する最適素片を選択する素片選択過程と、選択した最適素片を所定順に接続して合成音を生成する過程とを有する規則音声合成方法において、

前記素片選択過程は、

合成パラメタとなる一の波形素片の後続音素が合成時の音韻環境と一致しているか否かを判定するステップと、後続音素一致と判定された波形素片に対し、合成対象語彙の発声環境と該波形素片の抽出環境との差分を表す抽出誤差を用いて少なくとも一つの最適素片候補を選択するステップと、

選択された最適素片候補の韻律情報と所定の選択基準値との差分を表す選択誤差に基づいて最適素片を決定するステップとを有することを特徴とする規則音声合成方法。

【請求項 2】 前記最適素片を決定するステップは、個々の最適素片候補について前記抽出誤差及び選択誤差を順位付けるとともに、前記抽出誤差の順位に所定の係数を乗じて得た第 1 の値と前記選択誤差の順位に所定の係数を乗じて得た第 2 の値との合算値が最小となる最適素片候補を前記最適素片として決定することを特徴とする請求項 1 記載の規則音声合成方法。

【請求項 3】 前記決定された最適素片の韻律情報と前記選択基準値とに基づいて素片変形率を導出するステップと、

導出された素片変形率と所定の変形率しきい値とを比較し、素片変形率が変形率しきい値以下の場合は選択した最適素片をそのまま接続し、素片変形率が変形率しきい値を越える場合は該最適素片についての前記素片変形率を変形率しきい値以下に変形するステップと、を更に有することを特徴とする請求項 1 又は 2 記載の規則音声合成方法。

【請求項 4】 音韻単位に属する最適素片を選択する素片選択手段を備え、選択した最適素片を所定順に接続して合成音を生成する音声合成装置であって、前記素片選択手段は、合成パラメタとなる波形素片の韻律情報及び接続対象となる後続音素情報を含む素片情報を音韻単位毎に格納した素片情報辞書と、各音韻単位の分布統計情報及び選択時の読み込み数を定める情報を格納した音韻単位情報テーブルと、合成時に前記素片情報辞書及び音韻単位情報テーブルから最適素片に関する情報を索出する素片選択部とを含み、該素片選択部は、一の波形素片の後続音素が合成時の音韻環境と一致しているか否かを判定し、後続音素一致と判定された波形素片に対し、合成対象語彙の発声環境と該波形素片の抽出環境との差分を表す抽出誤差を用いて少なくとも一つの最適素片候補を選択するとともに、選択された最適素片候補の韻律情報と所定の選択基準値との差分を表す選択誤差に基づいて最適素片を決定することを特徴とする音

声合成装置。

【請求項 5】 前記素片選択手段で選択した最適素片の接続に先立ち、個々の最適素片の韻律情報を変形させて前記選択基準値との差を零値に近づける素片変形部と、前記素片選択手段で選択した最適素片の韻律情報と前記選択基準値とに基づいて素片変形率を導出する手段、及び導出された素片変形率が所定の変形率しきい値よりも大きいときのみ該選択された最適素片を前記素片変形部へ導く手段を備えた素片接続判定部と、を有することを特徴とする請求項 4 記載の音声合成装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、任意の入力文字列から合成音声生成する規則音声合成技術に関し、特に、合成パラメタである波形素片の選択手法及び選択した波形素片の変形手法に関する。

【0002】

【従来の技術】音声合成技術は、例えば、駅構内でのアナウンスや機械による文章朗読等に広く用いられている。音声合成に際しては、合成される音声を明瞭かつ違和感の無いものとするのが求められている。

【0003】従来の一般的な音声合成装置の機能ブロックの構成例を図 9 に示す。図 9 の構成の音声合成装置において、入力文字列は、入力端子 91 から入力された後に、前処理部 92 で日本語辞書 981 を参照して音韻単位に分割される。韻律設定部 93 は、その音韻単位の韻律パラメタを設定する。素片選択部 94 は、素片辞書 982 を参照して、上記設定された韻律パラメタの基準値に最も近い波形素片を選択する。素片変形部 95 においては、無条件で、選択された素片が基準値に合致するように、文章及び単語テキストの音声データを格納している素片ファイル 983 からの素片を夫々変形処理する。素片接続部 96 では、変形された素片をそれぞれ入力文字列の順に接続し、出力端子 97 を通じて出力する。これにより入力文字列に対応する合成音声を得られる。

【0004】

【発明が解決しようとする課題】上述した従来の音声合成装置において、素片選択部 94 では、合成時の音韻環境のような要因は考慮されておらず、その選択結果は、設定された韻律パラメタの基準値のみに強く依存している。そのため、語尾に近い音韻単位に対応する波形素片の選択時に、語頭に近い波形素片を選択してしまう現象が生じ、合成音声が不自然になる場合が多いという問題があった。また、従来、素片変形部 95 では、波形素片の選択結果とは無関係に、先に選択された基準値に合致するよう変形処理を行っているため、合成される音声の品質が劣化することもあった。

【0005】本発明の課題は、最適な波形素片を選択して合成時の自然性の劣化を防ぐことができる規則音声合成方法を提供することにある。本発明の他の課題は、波

形素片の過度の変形による品質劣化を抑え、従来装置に比べ、より自然性の高い合成音声を得られる構成の音声合成装置を提供することにある。

【0006】

【課題を解決するための手段】上記課題を解決する本発明の規則音声合成方法は、音韻単位に属する最適素片を選択する素片選択過程と、選択した最適素片を所定順に接続して合成音を生成する過程とを有する方法において、前記素片選択過程が、合成パラメタとなる一の波形素片の後続音素が合成時の音韻環境と一致しているか否かを判定するステップと、後続音素一致と判定された波形素片に対し、合成対象語彙の発声環境と該波形素片の抽出環境との差分を表す抽出誤差を用いて少なくとも一つの最適素片候補を選択するステップと、選択された最適素片候補の韻律情報と所定の選択基準値との差分を表す選択誤差に基づいて最適素片を決定するステップとを有することを特徴とする。

【0007】前記最適素片を決定するステップのより具体的な態様としては、例えば、個々の最適素片候補について前記抽出誤差及び選択誤差を順位付けるとともに、前記抽出誤差の順位に所定の係数を乗じて得た第1の値と前記選択誤差の順位に所定の係数を乗じて得た第2の値との合算値が最小となる最適素片候補を前記最適素片として決定する。

【0008】本発明の規則音声合成方法は、更に、前記決定された最適素片の韻律情報と前記選択基準値とに基づいて素片変形率を導出するステップと、導出された素片変形率と所定の変形率しきい値とを比較し、素片変形率が変形率しきい値以下の場合は選択した最適素片をそのまま接続し、素片変形率が変形率しきい値を超える場合は該最適素片についての前記素片変形率を変形率しきい値以下に変形するステップと、有することを特徴とする。

【0009】また、上記課題を解決する本発明の音声合成装置は、音韻単位に属する最適素片を選択する素片選択手段を備え、選択した最適素片を所定順に接続して合成音を生成する装置であって、前記素片選択手段が、合成パラメタとなる波形素片の韻律情報及び接続対象となる後続音素情報を含む素片情報を音韻単位毎に格納した素片情報辞書と、各音韻単位の分布統計情報及び選択時の読み込み数を定める情報を格納した音韻単位情報テーブルと、合成時に前記素片情報辞書及び音韻単位情報テーブルから最適素片に関する情報を索出する素片選択部とを含み、該素片選択部は、一の波形素片の後続音素が合成時の音韻環境と一致しているか否かを判定し、後続音素一致と判定された波形素片に対し、合成対象語彙の発声環境と該波形素片の抽出環境との差分を表す抽出誤差を用いて少なくとも一つの最適素片候補を選択するとともに、選択された最適素片候補の韻律情報と所定の選択基準値との差分を表す選択誤差に基づいて最適素片を

決定することを特徴とする。ここで、音韻単位情報テーブルに格納する分布統計情報とは、例えば個々の音韻単位についての韻律情報の分布を統計的に表したもので、上記選択誤差を導出する際に用いられるものである。また、選択時の読み込み数を定める情報とは、個々の波形素片について選択が予定される候補数を指標する情報である。

【0010】このような構成の音声合成装置によれば、音韻単位に属する最適素片の選択基準として、韻律情報だけでなく、もとの発話単位（又は語彙）の長さや波形素片（又は音韻）の占める位置のような、音韻単位や波形素片の抽出時の環境が加味されるので、合成時の不自然性が解消される。

【0011】本発明の音声合成装置は、また、前記素片選択手段で選択した最適素片の接続に先立ち、個々の最適素片の韻律情報を変形させて前記選択基準値との差を零値に近づける素片変形部と、前記素片選択手段で選択した最適素片の韻律情報と前記選択基準値とに基づいて素片変形率を導出する手段、及び導出された素片変形率が所定の変形率しきい値よりも大きいときのみ該選択された最適素片を前記素片変形部へ導く手段を備えた素片接続判定部と、を有するものである。このように構成することで、波形素片を変形する際の変形の程度が小さいものに対しては、最適素片が変形されずにそのまま用いられるので、過度の素片変形による音声の自然性劣化が防止される。

【0012】

【発明の実施の形態】以下、図面を参照して本発明の実施の形態を詳細に説明する。図1は、本発明の一実施形態を示す音声合成装置のブロック構成図である。なお、前処理部12、韻律設定部13、日本語辞書191の機能は、基本的には図9に示した従来装置のものと同様であり、その詳細な説明は省略する。

【0013】また、図2及び図3は、この音声合成装置における素片辞書192の詳細な構成例、図4～図7は、素片選択部14の動作原理を示すフローチャート、図8は、接続条件判定部25の動作原理を示すフローチャートである。以下、これらの図を用いて、本実施形態による音声合成処理の概要を説明する。

【0014】図1において、例えば読み上げ対象となる文字列は、入力端子11から入力される。前処理部12においては、日本語辞書191を用いて、入力された文字列（入力文字列）を解析し、合成単位である音韻単位及びその音韻単位に関するアクセント情報などを出力する。なお、本実施形態により使用される日本語辞書は、単語単位の読み及びアクセント型が記述された、音声合成用にカスタマイズされた辞書である。この辞書構成については、従来から種々のものが提案されている。また、本実施形態における音韻単位は、好ましくは、/a/や/k a/などの音節の他に、連母音/a i/や複合

音節／aN／などを含む。

【0015】次に、韻律設定部13では、音韻単位の韻律情報、例えば韻律パラメタの基準値を設定する。韻律情報（以下、韻律パラメタとする）は、素片選択部14で用いられる。韻律パラメタの要素としては、平均ピッチ周波数、ピッチ傾斜、時間長、平均パワーの4種類が用いられるが、この4種類の組み合わせによる新しいパラメタの算出、また4種類の取捨選択などは、合成音声の韻律条件に応じて適宜変更可能である。

【0016】素片選択部14では、韻律設定部13で設定された基準値と波形素片の情報を格納した素片辞書192とを用いて、音韻単位の最適素片を選択する。素片選択部14の詳細動作を、図2～図7を用いて説明する。

【0017】本実施形態で使用する素片辞書192は、図2に示したように、合成パラメタとなる波形素片の個々の情報を格納した素片情報辞書22と、各音韻単位の統計情報を記録した音韻単位情報テーブル23とを有し、これらが素片選択部14につながる入力部21及び出力部24に対して各々並列に接続されている。素片情報辞書22は、具体的には図3(a)のように構成され、音韻単位31、ファイル番号32、後続音素33、複数の辞書要素34、35を含んでいる。辞書要素34は、例えば平均ピッチ、ピッチ傾斜、RMSパワー、開始サンプル、終了サンプルから成り、辞書要素35は、例えば発話単位長及び発話単位位置より成る。

【0018】音韻単位情報テーブルは図3(b)のように構成され、音韻単位31、開始インデックス36、終了インデックス37、及び音韻単位要素38を含んでいる。開始インデックス36と終了インデックス37は、音韻単位に対応する波形素片の読み込み数を規定するものである。音韻単位要素38は、個々の音韻単位について*

抽出誤差＝

重み係数w1×発話単位内モーラ差分+重み係数w2×単位内位置差分…(1)

)

発話単位内モーラ差分＝

発話単位長－音韻単位数…(2)

単位内位置差分＝

発話単位位置－音韻単位位置…(3)

【0022】ここで、発話単位長及び発話単位位置とは、波形素片の抽出環境に関するもので、波形素片がどの程度の長さの発話単位のどの位置から抽出されたかを示す変数である。発話単位とは、一つの息継ぎの間に発声される音声情報を表す単位である。同様に、音韻単位数及び音韻単位位置は、選択対象となっている該音韻単位が、どの程度の長さの合成語彙のどの位置に属するかを示す変数である。抽出誤差算出の一例を図5に示す。

【0023】図5は、合成の対象となる単語“温かみ(あたかみ)”の2番目の音韻単位51(／ta／)

*での分布統計情報を規定するもので、例えば平均ピッチ、ピッチ傾斜、時間長、RMSパワーの最大値と最小値より成る。これらの最大値及び最小値は固定値とし、この範囲での情報を選択誤差の演算に用いる。なお、図3において“B”はバイトを示すが、図示のバイト数は例示であって、これに限定する趣旨ではない。

【0019】以下、素片選択部14の動作原理を、図4～図7を参照して説明する。図4を参照すると、素片選択部14は、該音韻単位に関する情報を、素片情報辞書22(図3の音韻単位31～辞書要素35に対応)及び音韻単位情報テーブル23(図3の音韻単位31、開始インデックス36～音韻単位要素38に対応)より読み込む(S(処理ステップ、以下同じ)101)。読み込み数は、音韻単位情報テーブル23の開始インデックス36及び終了インデックス37より算出される。

【0020】次に、当該音韻単位に属する波形素片の中から、最適素片となる候補を選出するとともに、複数の選択基準の1つである抽出誤差を算出する(S102～S104)。具体的には、音韻環境に関する適正チェックとして、該波形素片の後続音素が合成時の音韻環境と一致しているか否かのチェックを行う(S102)。この場合、合成時の後続音素と調音様式の同一である音素も一致とみなす。例えば合成時の音韻環境が／k／であった場合、調音様式の同一である／t／や／p／も一致とみなす。次に、S102において後続音素一致と判定された該波形素片に対して、抽出誤差を算出する(S103)。ここで、抽出誤差とは、合成対象となる語彙の発声環境と、波形素片の抽出環境との距離尺度となるものであり、この実施形態では、下記式(1)～(3)のように算出した。

【0021】

【数1】

を選択する際の抽出誤差を算出する方法を説明したものである。図5において、境界線を引かれている単位が一つの音韻単位に対応する。この場合、波形素片52に関して、発話単位長は“8”、発話単位位置は“4”となり、音韻単位51に関して、音韻単位数は“5”、音韻単位位置は“2”となるので、式(2)、(3)より、発話単位内モーラ差分は“8”－“5”＝“3”、単位内位置差分は“4”－“2”＝“2”となり、該波形素片52に関する抽出誤差は、(1)式より、w1×3＋w2×2となる。

【0024】図4に戻り、上記S102～S103の処理

を、読み込んだ素片の数すべてに対して繰り返す（S104）。その後、S104を経て選出された最適素片候補に対して、基準値との距離尺度である選択誤差を算出し（S105）、抽出誤差及び選択誤差から最適素片を決定する（S106）。選択誤差は、基準値及び音韻単

選択誤差

= 平均ピッチ誤差 + ピッチ傾斜誤差 + 時間長誤差 + RMSパワー誤差… (4)

各パラメータ誤差

= (正規化誤差値) × 重み係数… (5)

平均ピッチに関する正規化誤差値

= (基準値 - 素片値) / (素片最大値 - 素片最小値) … (6)

【0026】(6)式において、正規化誤差値の算出時に分母が“0”となる場合（該音韻単位に属する波形素片が1つしかない場合など）には、正規化誤差値を“0”とする。なお、選択誤差の算出方法には特に制限はなく、上記(4)～(6)式の方法以外にも、種々の方法を用いてよい。

【0027】S106の動作の詳細を、図6及び図7を参照して説明する。図6は、上記S106における処理手順を示すフローチャートであり、まず、候補番号の各※20

結合スコア

= 選択誤差順位 × 重み係数w1' + 抽出誤差順位 × 重み係数w2' … (7)

【0029】更に、最適素片候補に対する結合スコアを最小順に順位付けし（S204）、結合スコアの順位が最小となる候補を、最適素片として決定する（S205）。

【0030】図7は、特定の音韻単位に対する順位付けの結果の例であり、上記S104の終了時点で、候補番号0～11で示される12個の最適素片候補が選出された場合の例が示されている。図示の例では、w1' = 2、w2' = 1という設定になっている。これらの最適素片候補について、以上述べたS101～S105、及びS201～S205に詳細に示したS106の処理をすべての音韻単位について繰り返すことで（S107）、図7に斜線で示した候補番号“1”のものが最適素片として決定される。これにより、語頭に近い音韻単★

素片変形率 = |(基準値 / 素片値) - 1| × 100… (8)

【0034】次に、素片変形率と変形率しきい値を、各韻律パラメータについて比較し（S304）、その比較結果に対して、変形フラグ設定処理を行う（S305 / S306）。変形率しきい値は、音韻単位毎に変更可能な数値である。また、この例では波形素片の変形手法は一種であるが、波形素片の変形率に応じて変形手法を変えることで、より劣化の少ない合成音声を生成することも可能である。以上述べたS301からS305又はS306の処理を、全ての音韻単位に対して行う（S307）。最後に、該最適素片を素片ファイル193より切り出し、合成音声バッファに書き込む（S308）。

【0035】各韻律パラメータに関して、変形フラグ値が“1”の場合は、素片変形部16において、最適素片

※位情報テーブル23内の最適素片候補の韻律パラメータを用いて、式(4)～(6)のように算出される。

【0025】

【数2】

※々に対して最適素片候補に対する抽出誤差の順位付けを最小順に行う（S201）。次に、選択誤差に対する順位付けを同様に行って選択誤差順位を決定する（S202）。その後、各誤差の順位に基づいて、最適性の基準となる結合スコアを以下の式(7)により算出する（S203）。

【0028】

【数3】

★位に対しては、確実に語頭に近い位置にある波形素片が選択されるようになる。

【0031】図1に戻り、接続条件判定部15では、上記素片選択部14での処理結果に基づいて最適素片を変形するか否かの判定を行い、判定結果に対応した変形フラグ値を設定する。接続条件判定部15の詳細動作手順を図8に示す。

【0032】図8を参照すると、接続条件判定部15では、韻律設定部13で設定された上記基準値を読み込み（S301）、さらに最適素片の韻律パラメータを読み込む（S302）。次に、式(8)により、素片変形率を各韻律パラメータについて算出する（S303）。

【0033】

【数4】

を式(8)で算出した素片変形率に従って変形する。変形処理は、任意の手法が考えられるが、一例として、時間長変形の場合には、変形率に従って波形素片の該波形サンプルを間引き／補間する処理があげられる。変形された波形サンプルは、波形変形バッファ（図示省略）に別途格納される。素片フラグ値が0の場合は、変形処理を行わない。

【0036】図1の素片接続部17では、図示しない合成音声バッファ及び波形変形バッファに格納された波形サンプルを順次結合し、合成音声を生成する。以上述べた処理を経て、出力端子18には、入力文字列に対応した合成音声出力される。

【0037】

【発明の効果】以上の説明から明らかなように、本発明によれば、波形素片や音韻単位の抽出環境を含めた複数の選択基準により最適素片が選択されるので、合成される音声の自然性が向上する効果がある。また、最適素片の接続に先立ち、最適素片の変形処理の要／不要が判定され、変形処理の必要性が高い場合にのみ変形処理がなされるので、過度の変形処理による合成音声の品質劣化が回避される効果がある。

【図面の簡単な説明】

【図1】本発明の一実施形態となる音声合成装置のブロック構成図。

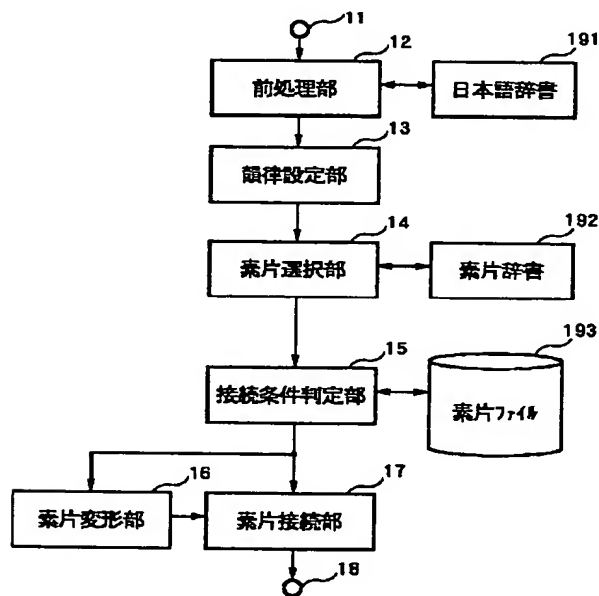
【図2】本実施形態の音声合成装置における素片辞書の構成例を示す説明図。

【図3】本実施形態の音声合成装置における素片辞書の内容説明図。

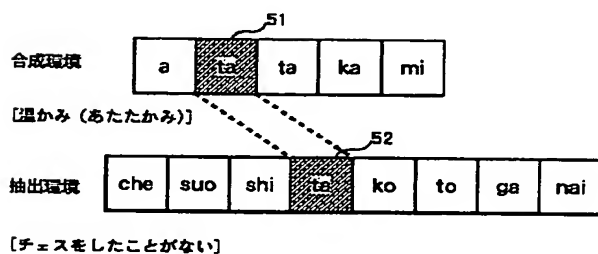
【図4】本実施形態の音声合成装置における素片選択部の動作原理を示すフローチャート。

【図5】抽出誤差の算出例の説明図。

【図1】



【図5】



【図6】最適素片の決定手順を示すフローチャート。

【図7】素片選択部における最適素片選択結果例の説明図。

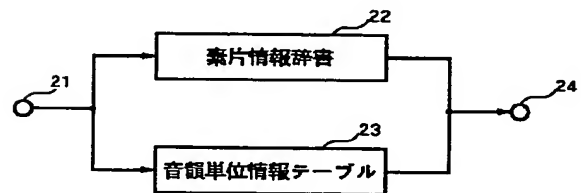
【図8】接続条件判定部の処理手順を示すフローチャート。

【図9】従来の音声合成装置のブロック構成図。

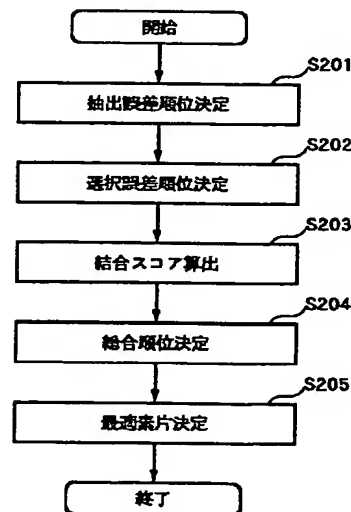
【符号の説明】

- 12 前処理部
- 13 韻律設定部
- 14 素片選択部
- 15 接続条件判定部
- 16 素片変形部
- 17 素片接続部
- 22 素片情報辞書
- 23 音韻単位情報テーブル
- 191 日本語辞書
- 192 素片辞書
- 193 素片ファイル

【図2】



【図6】



【図3】

(a) 素片情報辞書

音韻 単位 (8B)	フリガ 番号 (10B)	後続 音素 (8B)	平均 E _{eff} (4B)	E _{eff} 傾斜 (2B)	RMS パワ (4B)	開始 インデ (4B)	終了 インデ (4B)	発話 単位長 (1B)	発話単 位位置 (1B)
#ba-	sa0001	k	250.0	0.3	3000	4500	6500	8	1

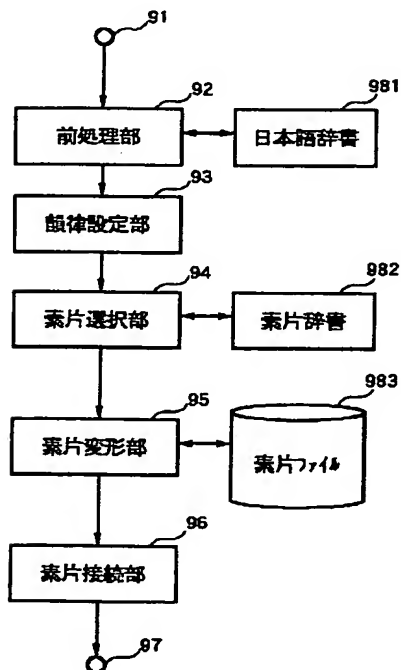
(b) 音韻単位情報テーブル

音韻 単位 (8B)	開始 インデックス (2B)	終了 インデックス (2B)	平均 E _{eff} 最大 最小 (2B) (2B)	E _{eff} 傾斜 最大 最小 (4B) (4B)	時間長 最大 最小 (2B) (2B)	RMS パワ 最大 最小 (2B) (2B)
#ba-	100	200	350 200	0.3 -0.3	250 200	2250 2000

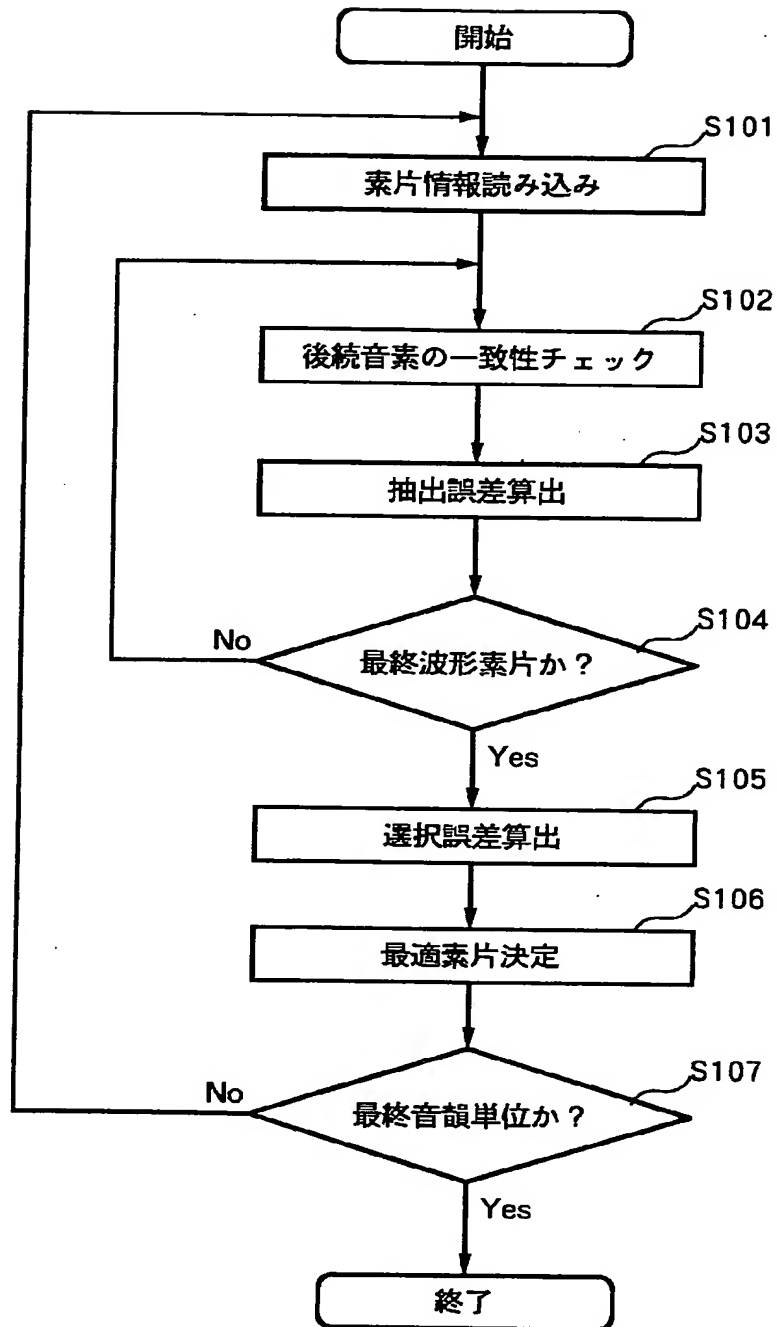
【図7】

候補番号	選択誤差順位	選択誤差	抽出誤差順位	抽出誤差	結合 ₁₂₇	総合順位
0	7	3.45	1	0	15	5
1	3	1.58	4	3	10	1
2	12	9.26	6	6	30	10
3	11	8.85	10	9	32	11
4	9	4.52	7	7	25	9
5	8	3.98	2	2	18	8
6	6	3.45	4	3	16	7
7	5	3.11	2	2	12	3
8	4	2.67	7	7	15	6
9	2	1.47	7	7	11	2
10	10	5.73	12	11	32	12
11	1	1.12	11	10	13	4

【図9】



【図4】



【図8】

